Assignment1 (Score: 9.0 / 9.0)

1. Test cell (Score: 1.0 / 1.0)
2. Test cell (Score: 1.0 / 1.0)
3. Test cell (Score: 1.0 / 1.0)
4. Test cell (Score: 1.0 / 1.0)
5. Test cell (Score: 1.0 / 1.0)
6. Test cell (Score: 1.0 / 1.0)
7. Test cell (Score: 1.0 / 1.0)
8. Test cell (Score: 1.0 / 1.0)
9. Test cell (Score: 1.0 / 1.0)

*You are currently looking at **version 0.1** of this notebook. To download notebooks and datafiles, as well as get help on Jupyter notebooks in the Coursera platform, visit the Jupyter Notebook FAQ course resource.*

# Assignment 1 - Introduction to Machine Learning ¶

For this assignment, you will be using the Breast Cancer Wisconsin (Diagnostic) Database to create a classifier that can help diagnose patients. First, read through the description of the dataset (below).

```
In [1]:  import numpy as np
         import pandas as pd
         from sklearn.datasets import load_breast_cancer

         cancer = load_breast_cancer()

         print(cancer.DESCR) # Print the data set description
```

```
.. _breast_cancer_dataset:

Breast cancer wisconsin (diagnostic) dataset
--------------------------------------------

**Data Set Characteristics:**

    :Number of Instances: 569

    :Number of Attributes: 30 numeric, predictive attributes and the class

    :Attribute Information:
        - radius (mean of distances from center to points on the perimeter)
        - texture (standard deviation of gray-scale values)
        - perimeter
        - area
        - smoothness (local variation in radius lengths)
        - compactness (perimeter^2 / area - 1.0)
        - concavity (severity of concave portions of the contour)
        - concave points (number of concave portions of the contour)
        - symmetry
        - fractal dimension ("coastline approximation" - 1)

        The mean, standard error, and "worst" or largest (mean of the three
        worst/largest values) of these features were computed for each image,
        resulting in 30 features.  For instance, field 0 is Mean Radius, field
        10 is Radius SE, field 20 is Worst Radius.

        - class:
                - WDBC-Malignant
                - WDBC-Benign

    :Summary Statistics:

    ===================================== ====== ======
                                           Min    Max
    ===================================== ====== ======
    radius (mean):                        6.981  28.11
    texture (mean):                       9.71   39.28
    perimeter (mean):                     43.79  188.5
    area (mean):                          143.5  2501.0
```

```
        smoothness (mean):                     0.053  0.163
        compactness (mean):                    0.019  0.345
        concavity (mean):                      0.0    0.427
        concave points (mean):                 0.0    0.201
        symmetry (mean):                       0.106  0.304
        fractal dimension (mean):              0.05   0.097
        radius (standard error):               0.112  2.873
        texture (standard error):              0.36   4.885
        perimeter (standard error):            0.757  21.98
        area (standard error):                 6.802  542.2
        smoothness (standard error):           0.002  0.031
        compactness (standard error):          0.002  0.135
        concavity (standard error):            0.0    0.396
        concave points (standard error):       0.0    0.053
        symmetry (standard error):             0.008  0.079
        fractal dimension (standard error):    0.001  0.03
        radius (worst):                        7.93   36.04
        texture (worst):                       12.02  49.54
        perimeter (worst):                     50.41  251.2
        area (worst):                          185.2  4254.0
        smoothness (worst):                    0.071  0.223
        compactness (worst):                   0.027  1.058
        concavity (worst):                     0.0    1.252
        concave points (worst):                0.0    0.291
        symmetry (worst):                      0.156  0.664
        fractal dimension (worst):             0.055  0.208
        ===================================== ====== ======

        :Missing Attribute Values: None

        :Class Distribution: 212 - Malignant, 357 - Benign

        :Creator:  Dr. William H. Wolberg, W. Nick Street, Olvi L. Mangasarian

        :Donor: Nick Street

        :Date: November, 1995

This is a copy of UCI ML Breast Cancer Wisconsin (Diagnostic) datasets.
https://goo.gl/U2Uwz2

Features are computed from a digitized image of a fine needle
```

aspirate (FNA) of a breast mass.  They describe
characteristics of the cell nuclei present in the image.

Separating plane described above was obtained using
Multisurface Method-Tree (MSM-T) [K. P. Bennett, "Decision Tree
Construction Via Linear Programming." Proceedings of the 4th
Midwest Artificial Intelligence and Cognitive Science Society,
pp. 97-101, 1992], a classification method which uses linear
programming to construct a decision tree.  Relevant features
were selected using an exhaustive search in the space of 1-4
features and 1-3 separating planes.

The actual linear program used to obtain the separating plane
in the 3-dimensional space is that described in:
[K. P. Bennett and O. L. Mangasarian: "Robust Linear
Programming Discrimination of Two Linearly Inseparable Sets",
Optimization Methods and Software 1, 1992, 23-34].

This database is also available through the UW CS ftp server:

ftp ftp.cs.wisc.edu
cd math-prog/cpo-dataset/machine-learn/WDBC/

.. topic:: References

    - W.N. Street, W.H. Wolberg and O.L. Mangasarian. Nuclear feature extraction
      for breast tumor diagnosis. IS&T/SPIE 1993 International Symposium on
      Electronic Imaging: Science and Technology, volume 1905, pages 861-870,
      San Jose, CA, 1993.
    - O.L. Mangasarian, W.N. Street and W.H. Wolberg. Breast cancer diagnosis and
      prognosis via linear programming. Operations Research, 43(4), pages 570-577,
      July-August 1995.
    - W.H. Wolberg, W.N. Street, and O.L. Mangasarian. Machine learning techniques
      to diagnose breast cancer from fine-needle aspirates. Cancer Letters 77 (1994)
      163-171.

The object returned by `load_breast_cancer()` is a scikit-learn Bunch object, which is similar to a dictionary.

```
In [2]: cancer.keys()

Out[2]: dict_keys(['data', 'target', 'frame', 'target_names', 'DESCR', 'feature_names', 'filename', 'data_mod
        ule'])
```

## Question 0 (Example)

How many features does the breast cancer dataset have?

*This function should return an integer.*

```
In [3]:  Student's answer                                                              (Top)

         # You should write your whole answer within the function provided. The autograder will call
         # this function and compare the return value against the correct solution value
         def answer_zero():
             # This function returns the number of features of the breast cancer dataset, which is an intege
         r.
             # The assignment question description will tell you the general format the autograder is expect
         ing

             return len(cancer['feature_names'])
             # YOUR CODE HERE
             #raise NotImplementedError()
         answer_zero()
         # You can examine what your function returns by calling it in the cell. If you have questions
         # about the assignment formats, check out the discussion forums for any FAQs
```

```
Out[3]: 30
```

```
In [4]:  Grade cell: cell-d2933751632e1611                                    Score: 1.0 / 1.0 (Top)
```

# Question 1

Scikit-learn works with lists, numpy arrays, scipy-sparse matrices, and pandas DataFrames, so converting the dataset to a DataFrame is not necessary for training this model. Using a DataFrame does however help make many things easier such as munging data, so let's practice creating a classifier with a pandas DataFrame.

Convert the sklearn.dataset `cancer` to a DataFrame.

*This function should return a `(569, 31)` DataFrame with*

*columns =*

```
['mean radius', 'mean texture', 'mean perimeter', 'mean area',
'mean smoothness', 'mean compactness', 'mean concavity',
'mean concave points', 'mean symmetry', 'mean fractal dimension',
'radius error', 'texture error', 'perimeter error', 'area error',
'smoothness error', 'compactness error', 'concavity error',
'concave points error', 'symmetry error', 'fractal dimension error',
'worst radius', 'worst texture', 'worst perimeter', 'worst area',
'worst smoothness', 'worst compactness', 'worst concavity',
'worst concave points', 'worst symmetry', 'worst fractal dimension',
'target']
```

*and index =*

```
RangeIndex(start=0, stop=569, step=1)
```

In [5]:

Student's answer                                                          (Top)

```python
def answer_one():
    # YOUR CODE HERE
    return pd.DataFrame(np.c_[cancer['data'], cancer['target']],
                    columns= np.append(cancer['feature_names'], ['target']))
answer_one()
    #raise NotImplementedError()
```

Out[5]:

| | mean radius | mean texture | mean perimeter | mean area | mean smoothness | mean compactness | mean concavity | mean concave points | mean symmetry | mean fractal dimension | ... | worst texture | per |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 17.99 | 10.38 | 122.80 | 1001.0 | 0.11840 | 0.27760 | 0.30010 | 0.14710 | 0.2419 | 0.07871 | ... | 17.33 | |
| 1 | 20.57 | 17.77 | 132.90 | 1326.0 | 0.08474 | 0.07864 | 0.08690 | 0.07017 | 0.1812 | 0.05667 | ... | 23.41 | |
| 2 | 19.69 | 21.25 | 130.00 | 1203.0 | 0.10960 | 0.15990 | 0.19740 | 0.12790 | 0.2069 | 0.05999 | ... | 25.53 | |
| 3 | 11.42 | 20.38 | 77.58 | 386.1 | 0.14250 | 0.28390 | 0.24140 | 0.10520 | 0.2597 | 0.09744 | ... | 26.50 | |
| 4 | 20.29 | 14.34 | 135.10 | 1297.0 | 0.10030 | 0.13280 | 0.19800 | 0.10430 | 0.1809 | 0.05883 | ... | 16.67 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 564 | 21.56 | 22.39 | 142.00 | 1479.0 | 0.11100 | 0.11590 | 0.24390 | 0.13890 | 0.1726 | 0.05623 | ... | 26.40 | |
| 565 | 20.13 | 28.25 | 131.20 | 1261.0 | 0.09780 | 0.10340 | 0.14400 | 0.09791 | 0.1752 | 0.05533 | ... | 38.25 | |
| 566 | 16.60 | 28.08 | 108.30 | 858.1 | 0.08455 | 0.10230 | 0.09251 | 0.05302 | 0.1590 | 0.05648 | ... | 34.12 | |
| 567 | 20.60 | 29.33 | 140.10 | 1265.0 | 0.11780 | 0.27700 | 0.35140 | 0.15200 | 0.2397 | 0.07016 | ... | 39.42 | |
| 568 | 7.76 | 24.54 | 47.92 | 181.0 | 0.05263 | 0.04362 | 0.00000 | 0.00000 | 0.1587 | 0.05884 | ... | 30.37 | |

569 rows × 31 columns

In [6]:

Grade cell: `cell-2dea923f2da8db76`                           Score: 1.0 / 1.0 (Top)

# Question 2

What is the class distribution? (i.e. how many instances of `malignant` and how many `benign` ?)

*This function should return a Series named* `target` *of length 2 with integer values and index =* `['malignant', 'benign']`

In [7]: 

Student's answer (Top)

```python
def answer_two():
    cancerdf = answer_one()
    malignant = len(cancerdf[cancerdf['target'] == 0])
    benign = len(cancerdf[cancerdf['target'] == 1])
    target = pd.Series(data = [malignant, benign], index = ['malignant', 'benign'])
    return target # Return your answer

answer_two()
    # YOUR CODE HERE
    #raise NotImplementedError()
```

Out[7]: 
```
malignant    212
benign       357
dtype: int64
```

In [8]: 

Grade cell: `cell-3d372226c8ec1345`                    Score: 1.0 / 1.0 (Top)

## Question 3

Split the DataFrame into  X  (the data) and  y  (the labels).

*This function should return a tuple of length 2:*  (X, y) *, where*

- X *has shape*  (569, 30)
- y *has shape*  (569,) .

In [9]:

Student's answer                                                                (Top)

```python
def answer_three():
    # YOUR CODE HERE
    cancerdf = answer_one()
    X = cancerdf.iloc[:, :30]
    y = cancerdf.iloc[:, 30]
    return X, y
    #raise NotImplementedError()
```

In [10]:

Grade cell: **cell-2ab04bcdf3007380**                          Score: 1.0 / 1.0 (Top)

## Question 4

Using `train_test_split`, split `X` and `y` into training and test sets `(X_train, X_test, y_train, and y_test)`.

**Set the random number generator state to 0 using `random_state=0` to make sure your results match the autograder!**

*This function should return a tuple of length 4:* `(X_train, X_test, y_train, y_test)`, *where*

- `X_train` *has shape* `(426, 30)`
- `X_test` *has shape* `(143, 30)`
- `y_train` *has shape* `(426,)`
- `y_test` *has shape* `(143,)`

In [11]: | Student's answer (Top)

```python
from sklearn.model_selection import train_test_split

def answer_four():
    # YOUR CODE HERE
    X, y = answer_three()

    # Your code here

    X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)
    return X_train, X_test, y_train, y_test
    #raise NotImplementedError()
```

In [12]: | Grade cell: `cell-725b24dae2118210`                           Score: 1.0 / 1.0 (Top)

## Question 5

Using KNeighborsClassifier, fit a k-nearest neighbors (knn) classifier with `X_train`, `y_train` and using one nearest neighbor ( `n_neighbors =` `1` ).

*This function should return a `sklearn.neighbors.classification.KNeighborsClassifier` .

In [13]: | Student's answer                                                                    (Top)

```python
from sklearn.neighbors import KNeighborsClassifier

def answer_five():
    # YOUR CODE HERE
    X_train, X_test, y_train, y_test = answer_four()

    knn = KNeighborsClassifier(n_neighbors = 1)
    return knn.fit(X_train, y_train)
    #raise NotImplementedError()
```

In [14]: | Grade cell: `cell-fe3813c4f3a2e07b`                        Score: 1.0 / 1.0 (Top)

In [15]: | answer_five()

Out[15]:
```
▼        KNeighborsClassifier
KNeighborsClassifier(n_neighbors=1)
```

# Question 6

Using your knn classifier, predict the class label using the mean value for each feature.

Hint: You can use `cancerdf.mean()[:-1].values.reshape(1, -1)` which gets the mean value for each feature, ignores the target column, and reshapes the data from 1 dimension to 2 (necessary for the precict method of KNeighborsClassifier).

In [16]:

```python
def answer_six():
    # YOUR CODE HERE
    cancerdf = answer_one()
    means = cancerdf.mean()[:-1].values.reshape(1, -1)
    predict = answer_five()
    label = predict.predict(means)
    return label
    # Return your answer
answer_six()
    #raise NotImplementedError()
```

Out[16]: array([1.])

In [17]:

# Question 7

Using your knn classifier, predict the class labels for the test set `X_test` .

*This function should return a numpy array with shape `(143,)` and values either `0.0` or `1.0` .*

In [18]: 

Student's answer

```python
def answer_seven():
    # YOUR CODE HERE
    X_train, X_test, y_train, y_test = answer_four()
    knn = answer_five()

    # Your code here

    return knn.predict(X_test)
answer_seven()

    #raise NotImplementedError()
```

Out[18]: array([1., 1., 1., 0., 1., 1., 1., 1., 1., 1., 0., 1., 1., 1., 0., 0., 1.,
       0., 0., 0., 0., 1., 1., 1., 0., 1., 1., 1., 1., 0., 1., 0., 1., 0.,
       1., 0., 1., 0., 1., 0., 0., 1., 0., 1., 0., 0., 1., 1., 1., 0., 0.,
       1., 0., 1., 1., 1., 1., 1., 1., 0., 0., 0., 1., 1., 0., 1., 0., 0.,
       0., 1., 1., 0., 1., 1., 0., 1., 1., 1., 1., 1., 0., 0., 0., 1., 0.,
       1., 1., 1., 0., 0., 1., 0., 1., 0., 1., 1., 0., 1., 1., 1., 1., 1.,
       1., 1., 0., 1., 0., 1., 0., 1., 1., 0., 0., 1., 1., 1., 0., 1., 1.,
       1., 1., 1., 1., 1., 0., 1., 1., 1., 1., 1., 0., 1., 1., 1., 1., 1.,
       1., 0., 0., 1., 1., 1., 0.])

In [19]: 
Grade cell: `cell-ece94681388729ef`                              Score: 1.0 / 1.0

## Question 8

Find the score (mean accuracy) of your knn classifier using `X_test` and `y_test` .

*This function should return a float between 0 and 1*

In [20]: 

```python
def answer_eight():
    # YOUR CODE HERE
    X_train, X_test, y_train, y_test = answer_four()
    knn = answer_five()

    # Your code here

    return knn.score(X_test, y_test)

answer_eight()
    #raise NotImplementedError()
```

Out[20]: 0.916083916083916

In [21]: 

## Optional plot

Try using the plotting function below to visualize the different predicition scores between train and test sets, as well as malignant and benign cells.

In [22]: Student's answer                                                   (Top)

```python
def accuracy_plot():
    import matplotlib.pyplot as plt

    %matplotlib notebook

    X_train, X_test, y_train, y_test = answer_four()

    # Find the training and testing accuracies by target value (i.e. malignant, benign)
    mal_train_X = X_train[y_train==0]
    mal_train_y = y_train[y_train==0]
    ben_train_X = X_train[y_train==1]
    ben_train_y = y_train[y_train==1]

    mal_test_X = X_test[y_test==0]
    mal_test_y = y_test[y_test==0]
    ben_test_X = X_test[y_test==1]
    ben_test_y = y_test[y_test==1]

    knn = answer_five()

    scores = [knn.score(mal_train_X, mal_train_y), knn.score(ben_train_X, ben_train_y),
              knn.score(mal_test_X, mal_test_y), knn.score(ben_test_X, ben_test_y)]


    plt.figure()

    # Plot the scores as a bar chart
    bars = plt.bar(np.arange(4), scores, color=['#4c72b0','#4c72b0','#55a868','#55a868'])

    # directly label the score onto the bars
    for bar in bars:
        height = bar.get_height()
        plt.gca().text(bar.get_x() + bar.get_width()/2, height*.90, '{0:.{1}f}'.format(height, 2),
                       ha='center', color='w', fontsize=11)

    # remove all the ticks (both axes), and tick labels on the Y axis
    plt.tick_params(top='off', bottom='off', left='off', right='off', labelleft='off', labelbottom
='on')

    # remove the frame of the chart
    for spine in plt.gca().spines.values():
```

```
        spine.set_visible(False)

    plt.xticks([0,1,2,3], ['Malignant\nTraining', 'Benign\nTraining', 'Malignant\nTest', 'Benign\nT
 est'], alpha=0.8);
    plt.title('Training and Test Accuracies for Malignant and Benign Cells', alpha=0.8)
```

In [23]:
```
# Uncomment the plotting function to see the visualization,
# Comment out the plotting function when submitting your notebook for grading

# accuracy_plot()
```

In [24]:
```
accuracy_plot()
```

In [ ]: